

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Using Packet Filters and Network Virtualization to
Restrict Network Communications**

Inventor(s):

**Aamer Hydrie
Galen C. Hunt
Steven P. Levi
Jakob Rehof
David S. Stutz
Bassam Tabbara
Mark D. VanAntwerp
Robert V. Welland**

ATTORNEY'S DOCKET NO. MS1-653US

004207-1-2356960

TECHNICAL FIELD

This invention relates to computer system management. More particularly, the invention relates to using packet filters and network virtualization to restrict network communications.

BACKGROUND OF THE INVENTION

Computers are becoming increasingly interconnected via a wide range of networks, including both public and private networks, such as local area networks (LANs), the Internet, etc. Although such interconnectivity can lead to a number of benefits, so too can it lead to problems. One predominant problem that can arise is that of security, such as how to keep a computer from accessing other computers it should not be accessing, how to keep other computers from accessing your computer, etc.

One specific area in which these security problems can arise is within "co-location facilities". A co-location facility refers to a complex that can house multiple servers, typically coupled to the Internet. The co-location facility typically provides a reliable Internet connection, a reliable power supply, and proper operating environment. The co-location facility also typically includes multiple secure areas (e.g., cages) into which different companies can situate their servers. The particular company is then responsible for managing the operation of the servers in their server cluster. These multiple servers can then operate together to make information available to client computers via the Internet. Security within such a co-location facility, however, is very important. For example, care should be taken to ensure that servers for one company housed at the facility cannot

1 communicate with servers for a competitor's company that are also housed at the
2 facility.

3 A "firewall" may be used to provide some security for computers.
4 However, firewalls typically operate to shield the outside world (e.g., the public
5 Internet) from the inside world (e.g., an internal private corporate LAN). Such
6 configurations thus do not prevent intra-LAN communications between different
7 computers within the corporate LAN.

8 Further firewalls (e.g., software firewalls) could also be installed at each
9 computer to provide security. However, current firewalls are typically designed to
10 prevent other computers from accessing the computer that they are installed on,
11 not restrict the computer's ability to access other computers. Some firewalls,
12 particularly those designed for home users, also employ parental controls.
13 Enabling parental controls allows a user of the computer (e.g., a parent) to restrict
14 the ability of that user or others (e.g., children) to access particular World Wide
15 Web sites on the Internet. However, such firewalls that are installed on a computer
16 are typically managed at the computer itself. Thus, the firewalls are susceptible to
17 being bypassed (or otherwise attacked) by a user of the computer. For example, a
18 user may erase or disable the firewall software, a user may load another operating
19 system that can bypass the firewall, etc. Thus, there still exists a need for
20 improved security among interconnected computers.

21 The invention described below addresses these disadvantages, using packet
22 filters and network virtualization to restrict network communications.
23
24
25

1 SUMMARY OF THE INVENTION

2 Using packet filters and network virtualization to restrict network
3 communications is described herein.

4 According to one aspect, a network mediator corresponding to a computing
5 device uses packet filters to restrict network communications. The network
6 mediator includes a set of one or more filters, each filter having parameters that
7 are compared to corresponding parameters of a data packet to be passed through
8 the network mediator (either from or to the computing device). The network
9 mediator determines whether to allow the data packet to pass through based on
10 whether the data packet parameters match any filter parameters. The set of filters
11 can be modified by a remote device, but cannot be modified by the computing
12 device whose communications are being restricted.

13 According to another aspect, a network mediator corresponding to a
14 computing device uses network virtualization to restrict network communications.
15 The network mediator maintains a mapping of virtual addresses to network
16 addresses, and allows the computing device to access only the virtual addresses.
17 When a data packet is sent from the computing device, the data packet includes
18 the virtual address which is then changed to the network address by the network
19 mediator prior to forwarding the packet on the network. Similarly, when a data
20 packet is received at the network mediator targeting the computing device, the
21 network mediator changes the network address in the data packet to the
22 corresponding virtual address. By virtualizing the addresses, the computing
23 device is restricted in its knowledge of and ability to access other devices over the
24 network because it has no knowledge of what the other devices' addresses are.
25

1 According to another aspect, a network mediator corresponding to a
2 computing device uses packet filters and a multiple managerial level architecture
3 to restrict network communications. The network mediator includes a set of one
4 or more filters, each filter having parameters that are compared to corresponding
5 parameters of a data packet to be passed through the network mediator (either
6 from or to the computing device). The network mediator then determines whether
7 to allow the data packet through based on whether the data packet parameters
8 match any filter parameters. The set of filters can be modified by remote devices
9 at different managerial levels. However, remote devices are prohibited from
10 modifying filters to make the filters less restrictive than filters imposed by higher
11 level devices.

12 13 **BRIEF DESCRIPTION OF THE DRAWINGS**

14 The present invention is illustrated by way of example and not limitation in
15 the figures of the accompanying drawings. The same numbers are used
16 throughout the figures to reference like components and/or features.

17 Fig. 1 shows a network environment such as may be used with certain
18 embodiments of the invention.

19 Fig. 2 is a block diagram illustrating a multiple-level filter administration
20 scheme in accordance with certain embodiments of the invention.

21 Fig. 3 is a block diagram illustrating an exemplary filter set and exemplary
22 virtualization data in accordance with certain embodiment of the invention.

23 Fig. 4 is a block diagram illustrating an exemplary co-location facility that
24 can include certain embodiments of the invention.
25

1 Fig. 5 is a block diagram illustrating an exemplary node of a co-location
2 facility in more detail in accordance with certain embodiments of the invention.

3 Fig. 6 is a flowchart illustrating an exemplary process for making
4 modifications to restrictions in a network mediator in accordance with certain
5 embodiments of the invention.

6 Fig. 7 is a flowchart illustrating an exemplary process for imposing
7 restrictions on communications of a computing device in accordance with certain
8 embodiments of the invention.

9 Fig. 8 shows a general example of a computer that can be used in
10 accordance with certain embodiments of the invention.

11 12 **DETAILED DESCRIPTION**

13 Fig. 1 shows a network environment such as may be used with certain
14 embodiments of the invention. The network environment in Fig. 1 is illustrated
15 with reference to two computing devices 102 and 104 communicating via a
16 network 106. Only two computing devices are illustrated in Fig. 1 so as not to
17 clutter the drawings. It is to be appreciated that the invention may be used with
18 any number of computing devices coupled together.

19 Computing devices 102 and 104 communicate with each other over a data
20 communications network 106. Communications network 106 can be any of a
21 variety of networks, including a public network (e.g., the Internet), as well as a
22 private network (e.g., a corporate local area network (LAN) or wide area network
23 (WAN)), or combinations thereof. Communications network 106 can be
24 implemented in any of a variety of different manners, including wired and/or
25 wireless communications media. Communications network 106 can be of any of a

1 wide variety of complexities, ranging from a single wire (e.g., plugged into a jack
2 on each of computing devices 102 and 104) to a complex network including
3 routers, bridges, both wired and wireless media, etc.

4 Communication over network 106 can be carried out using any of a wide
5 variety of communications protocols. In one implementation, computing devices
6 102 and 104 can communicate with one another using the Hypertext Transfer
7 Protocol (HTTP), in which World Wide Web pages are hosted by the devices 102
8 and 104 and written in a markup language, such as the Hypertext Markup
9 Language (HTML) or the eXtensible Markup Language (XML). The World Wide
10 Web (also referred to as simply the "web") is a collection of documents (also
11 referred to as "web pages") that users can view or otherwise render and which
12 typically include links to one or more other pages that the user can access.

13 Each computing device has a corresponding network mediator that can
14 restrict the ability of the device to communicate with other devices. The network
15 mediator may be a separate component (e.g., a router) coupled to the computing
16 device (e.g., network mediator 108 coupled to computing device 102) or
17 alternatively included as part of the computing device (e.g., network mediator 110
18 included in computing device 104). Network mediators 108 and 110 can be
19 implemented in any of a variety of manners, including software, firmware,
20 hardware, or combinations thereof.

21 Network mediator 110 can be implemented within computing device 104 in
22 any of a variety of manners. By way of example, network mediator 110 may be
23 implemented on a network interface card (NIC) of device 104, thereby allowing
24 filtering to occur at the point where network communications flow into and out of
25 device 104. By way of another example, computing device 104 may include a

processor(s) that supports multiple privilege levels (e.g., rings in an x86 architecture processor). The multiple rings provide a set of prioritized levels that software can execute at, often including 4 levels (Rings 0, 1, 2, and 3). Ring 0 is typically referred to as the most privileged ring. Software processes executing in Ring 0 can typically access more features (e.g., instructions) than processes executing in less privileged Rings. Furthermore, a processor executing in a particular Ring cannot alter code or data in a higher priority ring. Thus, network mediator 110 can be implemented to execute in Ring 0, while other software applications (including the operating system) execute in lower priority rings (e.g., Rings 1, 2, and/or 3). Thus, network mediator 110 is able to shield itself from other software applications, preventing those applications from modifying mediator 110 (e.g., by a rogue or malicious program trying to subvert the restrictions imposed by network mediator 110).

Each network mediator 108 and 110 includes a controller 112, a set of one or more filters 114, and optionally virtualization data 116. Controller 112 is the access point for the network mediator. Data packets desiring to be sent from or received at the corresponding computing device pass through controller 112, as do any requests to modify filters 114 or virtualization data 116.

Filters 114 are a set of one or more filters that impose restrictions on the ability of the corresponding computing device to transmit data packets to and/or receive data packets from other computing devices. Upon receipt of a data packet, controller 112 accesses filters 114 to determine whether one or more of filters 114 indicate that the data packet cannot be sent to (if the corresponding computing device is attempting to send the data packet) the targeted device or received from (if the corresponding computing device is the targeted device of the data packet)

1 the source device. If the data packet cannot be sent to the targeted device (or
2 received from the source device), then controller 112 refuses to let the data packet
3 through and drops the data packet. Alternatively, a message may be returned to
4 the source of the data packet informing the source that it cannot send the desired
5 packet.

6 Virtualization data 116 includes the data to map virtual addresses to
7 network addresses. Although illustrated as separate from filters 114, virtualization
8 data 116 may alternatively be incorporated into filters 114. The network address
9 of a device on network 106 uniquely identifies the device (or group of devices) on
10 network 106 and can be used by other devices to send data packets to that
11 device(s). The format of the network address can vary based on the protocol(s)
12 used by network 106 (e.g., if the Internet Protocol (IP) is supported by network
13 106, then the network addresses can be 32-bit IP addresses).

14 In embodiments using virtual addresses, a computing device uses virtual
15 addresses to identify other computing devices. These virtual addresses uniquely
16 identify other devices within a particular computing device, but can have no
17 relationship to the actual network address for those other devices. The virtual
18 addresses thus are relatively useless outside of the computing device (and network
19 mediator). When a data packet is to be sent to another computing device, network
20 mediator uses virtualization data 116 to map the virtual address to the correct
21 network address for the targeted computing device so that the data packet can be
22 communicated to the targeted computing device. Similarly, when a data packet
23 received from another computing device targeting the computing device
24 corresponding to the network mediator, the network mediator uses virtualization
25

1 data 116 to map the network address to the correct virtual address for the source
2 computing device.

3 The virtual addresses used by a computing device can be generated in any
4 of a variety of manners. By way of example, they may be generated randomly or
5 in accordance with some predetermined (or dynamic) algorithm. The virtual
6 addresses for a computing device can be generated by the network mediator
7 corresponding to that device, or alternatively some other device or component
8 coupled to the network mediator (e.g., via network 106).

9 Controller 112 restricts filters 114 and virtualization data 116 to being
10 managed only from one or more remote devices. Controller 112 prevents the
11 computing device that the network mediator corresponds to from modifying the
12 filters 114 and virtualization data 116 (e.g., prevents the computing device from
13 adding filters or virtualization data, removing filters or virtualization data,
14 modifying filters or virtualization data, etc.). Thus, the computing device that is
15 having its communication restricted by the network mediator does not have the
16 authority to alter any of the filters 114 or virtualization data 116 being used to
17 impose those restrictions.

18 Restricting managerial control to only remote devices can be implemented
19 in a variety of different manners. In one implementation, the network mediator
20 does not expose an interface to the corresponding computing device to allow the
21 computing device to make any requests to change the filters 114 or virtualization
22 data 116. The network mediator allows only authorized remote devices (e.g., only
23 devices with particular addresses and/or that can authenticate themselves using an
24 identifier and a password) to modify filters 114 and data 116. The network
25 mediator can be initialized to accept management commands from only a

1 particular authorized device (or user id), which in turn can configure network
2 mediator to authorize additional devices (or user id's) to manage the network
3 mediator. Furthermore, the filters 114 and data 116 can be stored by network
4 mediator in a protected area (e.g., memory or registers that are encrypted and/or
5 require device (or user) verification to access it). Such storage in a protected area
6 prevents any device from bypassing any authentication procedures imposed by the
7 network mediator.

8 Controller 112 may also optionally maintain a log or other record (not
9 shown) of data packets that were not allowed to pass through the network
10 mediator. This record may be a copy of the entire data packet, or alternatively
11 only selected parts of the data packet (e.g., source and/or destination address, data
12 packet protocol, etc.). Additional information may also be maintained, such as a
13 timestamp indicating the date and time the data packet was received, which filter
14 was the cause of the refusal to allow the packet through, etc. Such information
15 can then be used for a variety of different manners. For example, the information
16 could be examined to try to identify if a malicious user or program is attempting to
17 break into a particular computing device(s), or to identify an improperly
18 functioning program that, due to an error in the program, is attempting to access
19 computing devices it should not be (e.g., during a debugging process), etc.

20 Controller 112 also supports a multiple-level restriction administration
21 scheme. In support of such a scheme, controller 112 allows different remote
22 devices to have different managerial levels, and prevents devices at lower-level
23 managerial levels from modifying filters 114 in a manner that could result in a
24 violation of a filter imposed by a higher-level managerial level (or changing of
25 virtualization data 116 established by a higher-level managerial level).

Fig. 2 is a block diagram illustrating a multiple-level filter administration scheme in accordance with certain embodiments of the invention. Multiple different computing devices 130, 132, 134, and 136 are illustrated, each with a corresponding network mediator (NM). Two different managerial levels are illustrated in Fig. 2, one being higher level with management being handled from administrator device 138 and another being lower level with management being handled from sub-administrator device 140. Although illustrated as two separate devices 138 and 140, multiple managerial levels may alternatively be implemented using a single device but different user identifications (with the network mediator verifying authenticity of the user rather than the device).

Administrator device 138, being the highest level managerial device, can modify filters 114 and virtualization data 116 for any of the network mediators corresponding to the computing devices 130 – 136. Sub-administrator device 140, however, is of a lower managerial level than administrator device 138, and can manage network mediators corresponding to devices 134 and 136, but not network mediators corresponding to devices 130 and 132. Thus, the network mediators for devices 134 and 136 are subject to management from multiple managerial levels, while the network mediators for devices 130 and 132 are subject to management from only a single managerial level. Sub-administrator device 140 can modify filters 114 and virtualization data 116 in network mediators corresponding to devices 134 and 136, but only so long as they do not conflict with (e.g., would result in a violation of) modifications previously made by administrator device 138. Any change by a lower managerial level device (or user) of a virtual address previously established (or modified) by a higher managerial level device (or user) would be such a violation. Two filters conflict with each other if the filter added

(or modified) by the higher managerial level device (or user) is more restrictive than the filter added (or modified) by the lower managerial level. A first filter that allows any type of access that is not allowed by a second filter is said to be less restrictive than the second filter. By way of example, a first filter may allow communications to be sent to a device at a target address only of a particular protocol, while a second filter may allow communications to be sent to the device at the target address using any protocol. In this example, the second filter would be less restrictive than the first filter. Note, however, that two filters are not in conflict with each other if the filter added (or modified) by the lower managerial level device (or user) is less restrictive than the filter added by the higher managerial level device (or user).

Administrator device 138 and sub-administrator device 140 are merely the devices via which management commands are issued to the computing devices 130 – 136. Management commands may also be issued from other devices, and multiple devices may exist that can issue management commands for the same management level (e.g., multiple administrator devices 138 operating at the highest managerial level may exist). Although not illustrated in Fig. 2, devices 138 and 140 may be computing devices analogous to devices 102 and 104 of Fig. 1, with corresponding network mediators themselves. By way of example, sub-administrator device 140 may have a network mediator that can be managed by administrator device 138 to restrict device 140 to communicating only with devices 134, 136, and 138 (and in this manner preventing sub-administrator device 140 from issuing management commands to devices 130 and 132).

In order to prevent lower-level restrictions from being imposed which would violate higher-level restrictions, each network mediator maintains an

identification of which level imposed each filter (or mapping). This identification may be direct (e.g., an indicator associated with each filter identifying a particular managerial level) or alternatively indirect (e.g., different data structures may be established for different managerial levels).

Situations can arise where a filter added by a higher-level device is in conflict with a filter previously added by a lower-level device. These situations can be resolved in different manners. In one implementation, the network mediator compares a new filter (newly added or newly changed) to all previously imposed filters. If the new filter conflicts with any previously imposed lower-level filter, then the previously imposed lower-level filter is deleted. In another implementation, the network mediator maintains filters from different managerial levels in different structures (e.g., tables). Thus, during the filtering process, the network mediator controller can compare the data packet to the filters on a table-by-table basis (e.g., so a restriction imposed by a higher managerial level will cause the data packet to be dropped before a table including lower managerial level filters is accessed (and so the less restrictive lower level filter would never get the chance to allow the data packet to pass through)).

Fig. 3 is a block diagram illustrating an example filter set and virtualization data in accordance with certain embodiment of the invention. The filter set and virtualization data 160 of Fig. 3 are discussed with additional reference to components in Fig. 1. Although the filter set and virtualization data 160 are illustrated in a table format for ease of illustration, it is to be appreciated that filters and virtualization data 160 can be maintained at a network mediator 108 or 110 using any of a wide variety of conventional data structures. Furthermore, the filter set and virtualization data 160 are illustrated with reference to Internet

Protocol (IP) data packet filtering. Alternatively, data packets can be filtered for different protocols, with the parameters of the filters varying based on the protocol (e.g., there would be no port parameters if the protocol did not support ports).

Filter set and virtualization data 160 includes multiple parameters or fields, one or more of which may be filled in for a particular filter. The fields include: a source address field 162, a destination address field 164, a source port field 166, a destination port field 168, a protocol field 170, and a mapping field 172. Each one of the fields 162 – 170 can identify a particular filter parameter for the corresponding filter. These filter parameters for a filter are then compared to the corresponding parameters of a data packet to determine whether the packet satisfies the filter. A data packet satisfies a filter if the filter parameters match (are the same as) the corresponding parameters of the data packet. Also, in the illustrated example filter set and virtualization data 160, each of the filters 174, 176, and 178 is a "permissive" filter – a packet received at the network mediator that satisfies one of these filters will be passed through to its destination (and will be dropped if it does not satisfy any of the filters). Alternatively, filter set and virtualization data 160 could include only "exclusionary" filters (that is, any packet received at the network mediator that satisfies one of the filters will be dropped; otherwise, the packet will be allowed through to its destination), or a combination of permissive and exclusionary filters.

In the illustrated example filter 174, any data packet received by the network mediator (targeting the corresponding computing device) from another computing device having a source address of "152.48.72.0", a source port of "2789", and using the UDP (User Datagram Protocol) protocol will be allowed through to the corresponding computing device. Additionally, mapping field 172

1 indicates to controller 112 to change the source address of "152.48.72.0" to
2 "143.62.79.83" in the data packet before forwarding the data packet on to the
3 computing device (or another portion of the computing device) for processing.

4 Similarly, filter 176 indicates that any data packet requested to be sent to a
5 target computing device by the computing device corresponding to the network
6 mediator can be sent if the target device address (the destination address) is
7 "173.42.68.200" and uses the TCP (Transmission Control Protocol) protocol. If
8 the data packet satisfies these parameters, then mapping field 172 indicates to
9 controller 112 to change the destination address in the data packet to
10 "143.62.79.82" prior to forwarding the data packet to the targeted device.

11 The parameters for fields 162 – 170 may also employ "wild cards", which
12 allow at least a portion of a parameter to match anything. Wild cards can be
13 implemented in any of a wide variety of manners, and in the illustrated example
14 are implemented using a value and corresponding mask. The value 180 and mask
15 182 for destination address field 164 of filter 178 are illustrated by way of
16 example in Fig. 3. The address stored in field 164 of filter 178 ("152.48.0.0") is
17 stored as a 32-bit value 180. Each of the 32 bits has a corresponding bit in the 32-
18 bit mask value 182. For each bit of value 180, the corresponding mask bit in mask
19 value 182 indicates whether that bit must match in the data packet to satisfy the
20 filter. For each bit in mask value 182 that is set (e.g., has a value of "1"), the
21 corresponding bit in value 180 must match in the data packet to satisfy the filter,
22 and for each bit in mask value 182 that is not set (e.g., has a value of "0"), the
23 corresponding bit in value 180 need not match in the data packet to satisfy the
24 filter. Thus, in the illustrated example, the first sixteen bits of value 180
25 (corresponding to "152.48") must match the corresponding field of the data packet

in order to satisfy the filter, although the remaining sixteen bits do not (therefore, the values "0.0" can basically be ignored). Thus, by way of example, if the data packet had a destination address of "152.48.137.72", then the data packet would satisfy the filter (assuming the data packet also had a destination port of 1312 and used the TCP protocol). However, if the data packet had a destination address of "148.39.152.48", then the data packet would not satisfy the filter (regardless of the values in its destination port and protocol fields).

Various modifications can be made to filter set and virtualization data 160. By way of example, mapping field 172 may be separate from the filter fields 162-170 (e.g., stored in a separate table or other data structure). By way of another example, an additional indication (not shown) may be included for one or more filters to indicate whether the filter corresponds to incoming data packets (those packets targeting the computing device corresponding to the network mediator) or outgoing data packets (those packets which are trying to be sent from the computing device corresponding to the network mediator). By way of another example, an additional indication (not shown) may be included for one or more filters to indicate a particular device (or managerial) level that implemented the filter (e.g., device 138 or 140 of Fig. 2).

Fig. 4 is a block diagram illustrating an exemplary co-location facility that can include certain embodiments of the invention. Co-location facility 208 is illustrated including multiple nodes (also referred to as server computers) 210. Each one of these nodes 210 can be a computing device 102 or 104 of Fig. 1, and includes a corresponding network mediator. Co-location facility 208 can include any number of nodes 210, and can easily include an amount of nodes numbering into the thousands.

The nodes 210 are grouped together in clusters, referred to as server clusters (or node clusters). For ease of explanation and to avoid cluttering the drawings, only a single cluster 212 is illustrated in Fig. 4. Each server cluster includes nodes 210 that correspond to a particular customer of co-location facility 104. The nodes 210 of a server cluster are physically isolated from the nodes 210 of other server clusters. This physical isolation can take different forms, such as separate locked cages or separate rooms at co-location facility 104. Physically isolating server clusters ensures customers of co-location facility 104 that only they can physically access their nodes (other customers cannot).

A landlord/tenant relationship (also referred to as a lessor/lessee relationship) can also be established based on the nodes 210. The owner (and/or operator) of co-location facility 104 owns (or otherwise has rights to) the individual nodes 210, and thus can be viewed as a "landlord". The customers of co-location facility 104 lease the nodes 210 from the landlord, and thus each can be viewed as a "tenant". The landlord is typically not concerned with what types of data or programs are being stored at the nodes 210 by the tenant, but does impose boundaries on the clusters that prevent nodes 210 from different clusters from communicating with one another, as discussed in more detail below.

Although physically isolated, nodes 210 of different clusters are often physically coupled to the same transport medium (or media) 211 that enables access to network connection(s) 216, and possibly application operations management console 242, discussed in more detail below. This transport medium can be wired or wireless.

As each node 210 can be coupled to a shared transport medium 211, each node 210 is configurable to restrict which other nodes 210 data packets can be sent

1 to or received from. Given that a number of different nodes 210 may be included
2 in a tenant's server cluster, the tenant may want to be able to pass data between
3 different nodes 210 within the cluster for processing, storage, etc. However, the
4 tenant will typically not want data to be passed to other nodes 210 that are not in
5 the server cluster. Configuring each node 210 in the cluster to restrict which other
6 nodes 210 data packets can be sent to or received from allows a boundary for the
7 server cluster to be established and enforced. Establishment and enforcement of
8 such server cluster boundaries prevents tenant data from being erroneously or
9 improperly forwarded to a node that is not part of the cluster.

10 These initial boundaries established by the landlord prevent communication
11 between nodes 210 of different tenants, thereby ensuring that each tenant's data
12 can be passed only to other nodes 210 of that tenant. The tenant itself may also
13 further define sub-boundaries within its cluster, establishing sub-clusters of nodes
14 210 that data cannot be communicated out of (or in to) either to or from other
15 nodes in the cluster. The tenant is able to add, modify, remove, etc. such sub-
16 cluster boundaries at will, but only within the boundaries defined by the landlord
17 (that is, the cluster boundaries). Thus, the tenant is not able to alter boundaries in
18 a manner that would allow communication to or from a node 210 to extend to
19 another node 210 that is not within the same cluster.

20 Co-location facility 104 supplies reliable power 214 and reliable network
21 connection(s) 216 to each of the nodes 210. Power 214 and network connection(s)
22 216 are shared by all of the nodes 210, although alternatively separate power 214
23 and network connection(s) 216 may be supplied to nodes 210 or groupings (e.g.,
24 clusters) of nodes. Any of a wide variety of conventional mechanisms for
25 supplying reliable power can be used to supply reliable power 214, such as power

received from a public utility company along with backup generators in the event of power failures, redundant generators, batteries, fuel cells, or other power storage mechanisms, etc. Similarly, any of a wide variety of conventional mechanisms for supplying a reliable network connection can be used to supply network connection(s) 216, such as redundant connection transport media, different types of connection media, different access points (e.g., different Internet access points, different Internet service providers (ISPs), etc.).

Management of each node 210 is carried out in a multiple-tiered manner. The multi-tiered architecture includes three tiers: a cluster operations management tier, an application operations management tier, and an application development tier. The cluster operations management tier is implemented locally at the same location as the node(s) being managed (e.g., at a co-location facility) and involves managing the hardware operations of the node(s). The cluster operations management tier is not concerned with what software components are executing on the nodes 210, but only with the continuing operation of the hardware of nodes 210 and establishing any boundaries between clusters of nodes.

The application operations management tier, on the other hand, is implemented at a remote location other than where the server(s) being managed are located (e.g., other than the co-location facility), but from a client computer that is still communicatively coupled to the server(s). The application operations management tier involves managing the software operations of the server(s) and defining sub-boundaries within server clusters. The client can be coupled to the server(s) in any of a variety of manners, such as via the Internet or via a dedicated (e.g., dial-up) connection. The client can be coupled continually to the server(s), or alternatively sporadically (e.g., only when needed for management purposes).

1 The application development tier is implemented on another client
2 computer at a location other than the server(s) (e.g., other than at the co-location
3 facility) and involves development of software components or engines for
4 execution on the server(s). Alternatively, current software on a node 210 at co-
5 location facility 208 could be accessed by a remote client to develop additional
6 software components or engines for the node. Although the client at which
7 application development tier is implemented is typically a different client than that
8 at which application operations management tier is implemented, these tiers could
9 be implemented (at least in part) on the same client.

10 Co-location facility 208 includes a cluster operations management console
11 for each server cluster. In the example of Fig. 4, cluster operations management
12 console 240 corresponds to cluster 212. Cluster operations management console
13 240 implements the cluster operations management tier for cluster 212 and is
14 responsible for managing the hardware operations of nodes 210 in cluster 212.
15 Cluster operations management console 240 monitors the hardware in cluster 212
16 and attempts to identify hardware failures. Any of a wide variety of hardware
17 failures can be monitored for, such as processor failures, bus failures, memory
18 failures, etc. Hardware operations can be monitored in any of a variety of
19 manners, such as cluster operations management console 240 sending test
20 messages or control signals to the nodes 210 that require the use of particular
21 hardware in order to respond (no response or an incorrect response indicates
22 failure), having messages or control signals that require the use of particular
23 hardware to generate periodically sent by nodes 210 to cluster operations
24 management console 240 (not receiving such a message or control signal within a
25 specified amount of time indicates failure), etc. Alternatively, cluster operations

1 management console 240 may make no attempt to identify what type of hardware
2 failure has occurred, but rather simply that a failure has occurred.

3 Once a hardware failure is detected, cluster operations management console
4 240 acts to correct the failure. The action taken by cluster operations management
5 console 240 can vary based on the hardware as well as the type of failure, and can
6 vary for different server clusters. The corrective action can be notification of an
7 administrator (e.g., a flashing light, an audio alarm, an electronic mail message,
8 calling a cell phone or pager, etc.), or an attempt to physically correct the problem
9 (e.g., reboot the node, activate another backup node to take its place, etc.).

10 Cluster operations management console 240 also establishes cluster
11 boundaries within co-location facility 208 by adding filters to the network
12 mediator corresponding to each node 210 that allows the node to communicate
13 only with other nodes in its cluster. The cluster boundaries established by console
14 240 prevent nodes 210 in one cluster (e.g., cluster 212) from communicating with
15 nodes in another cluster (e.g., any node not in cluster 212), while at the same time
16 not interfering with the ability of nodes 210 within a cluster from communicating
17 with other nodes within that cluster. These boundaries provide security for the
18 tenants' data, allowing them to know that their data cannot be communicated to
19 other tenants' nodes 210 at facility 104 even though network connection 216 may
20 be shared by the tenants.

21 In the illustrated example, each cluster of co-location facility 104 includes a
22 dedicated cluster operations management console. Alternatively, a single cluster
23 operations management console may correspond to, and manage hardware
24 operations of, multiple server clusters. According to another alternative, multiple
25 cluster operations management consoles may correspond to, and manage hardware

operations of, a single server cluster. Such multiple consoles can manage a single server cluster in a shared manner, or one console may operate as a backup for another console (e.g., providing increased reliability through redundancy, to allow for maintenance, etc.).

An application operations management console 242 is also communicatively coupled to co-location facility 208. Application operations management console 242 is located at a location remote from co-location facility 208 (that is, not within co-location facility 208), typically being located at the offices of the customer. A different application operations management console 242 corresponds to each server cluster of co-location facility 208, although alternatively multiple consoles 242 may correspond to a single server cluster, or a single console 242 may correspond to multiple server clusters. Application operations management console 240 implements the application operations management tier for cluster 212 and is responsible for managing the software operations of nodes 210 in cluster 212 as well as securing sub-boundaries within cluster 212. Application operations management console 240 can create, modify, and remove sub-boundaries by modifying the filter set in the network mediator corresponding to each node in the cluster to allow communication only with the desired nodes.

Application operations management console 242 monitors the software in cluster 212 and attempts to identify software failures. Any of a wide variety of software failures can be monitored for, such as application processes or threads that are "hung" or otherwise non-responsive, an error in execution of application processes or threads, etc. Software operations can be monitored in any of a variety of manners (similar to the monitoring of hardware operations discussed above),

1 such as application operations management console 242 sending test messages or
2 control signals to particular processes or threads executing on the nodes 210 that
3 require the use of particular routines in order to respond (no response or an
4 incorrect response indicates failure), having messages or control signals that
5 require the use of particular software routines to generate periodically sent by
6 processes or threads executing on nodes 210 to application operations
7 management console 242 (not receiving such a message or control signal within a
8 specified amount of time indicates failure), etc. Alternatively, application
9 operations management console 242 may make no attempt to identify what type of
10 software failure has occurred, but rather simply that a failure has occurred.

11 Once a software failure is detected, application operations management
12 console 242 acts to correct the failure. The action taken by application operations
13 management console 242 can vary based on the hardware as well as the type of
14 failure, and can vary for different server clusters. The corrective action can be
15 notification of an administrator (e.g., a flashing light, an audio alarm, an electronic
16 mail message, calling a cell phone or pager, etc.), or an attempt to correct the
17 problem (e.g., reboot the node, re-load the software component or engine image,
18 terminate and re-execute the process, etc.).

19 Thus, the management of a node 210 is distributed across multiple
20 managers, regardless of the number of other nodes (if any) situated at the same
21 location as the node 210. The multi-tiered management allows the hardware
22 operations management to be separated from the application operations
23 management, allowing two different consoles (each under the control of a different
24 entity) to share the management responsibility for the node.
25

Fig. 5 is a block diagram illustrating an exemplary node of a co-location facility in more detail in accordance with certain embodiments of the invention. Node 210 includes a monitor 250, referred to as the "BMonitor", and a plurality of software components or engines 252, and is coupled to (or alternatively incorporates) a mass storage device 262. In the illustrated example of Fig. 5, BMonitor 250 acts as network mediator 108 or 110 of Fig. 1.

In Fig. 5, node 210 is a server computer having a processor(s) that supports multiple privilege levels (e.g., rings in an x86 architecture processor). In the illustrated example, these privilege levels are referred to as rings, although alternate implementations using different processor architectures may use different nomenclature. The multiple rings provide a set of prioritized levels that software can execute at, often including 4 levels (Rings 0, 1, 2, and 3), with Ring 0 being the most privileged ring. In the illustrated example, BMonitor 250 executes in Ring 0, while engines 252 execute in Ring 1 (or alternatively Rings 2 and/or 3). Thus, the code or data of BMonitor 250 (executing in Ring 0) cannot be altered directly by engines 252 (executing in Ring 1). Rather, any such alterations would have to be made by an engine 252 requesting BMonitor 250 to make the alteration (e.g., by sending a message to BMonitor 250, invoking a function of BMonitor 250, etc.). Implementing BMonitor 250 in Ring 0 protects BMonitor 250 from a rogue or malicious engine 252 that tries to bypass any restrictions imposed by BMonitor 250.

BMonitor 250 is the fundamental control module of node 210 – it controls (and optionally includes) both the network interface card and the memory manager. By controlling the network interface card (which may be separate from BMonitor 250, or alternatively BMonitor 250 may be incorporated on the network

1 interface card), BMonitor 250 can control data received by and sent by node 210.
 2 By controlling the memory manager, BMonitor 250 controls the allocation of
 3 memory to engines 252 executing in node 210 and thus can assist in preventing
 4 rogue or malicious engines from interfering with the operation of BMonitor 250.

5 Although various aspects of node 210 may be under control of BMonitor
 6 250 (e.g., the network interface card), BMonitor 250 still makes at least part of
 7 such functionality available to engines 252 executing on the node 210. BMonitor
 8 250 provides an interface (e.g., via controller 254 discussed in more detail below)
 9 via which engines 252 can request access to the functionality, such as to send data
 10 out to another node 210 or to the Internet. These requests can take any of a variety
 11 of forms, such as sending messages, calling a function, etc.

12 BMonitor 250 includes controller 254 (performing the functions of
 13 controller 112 of Fig. 1), one or more filters 114, virtualization data 116, network
 14 interface 256, one or more keys 258, and a Distributed Host Control Protocol
 15 (DHCP) module 260. Network interface 256 provides the interface between node
 16 210 and the network (e.g., network connections 216 of Fig. 4) via the internal
 17 transport medium 211 of co-location facility 104. Filters 114 identify other nodes
 18 210 and possibly other sources or targets (e.g., coupled to network 106 of Fig. 1)
 19 that data can (or alternatively cannot) be sent to and/or received from. The nodes
 20 or other sources/targets can be identified in any of a wide variety of manners, such
 21 as by network address (e.g., Internet Protocol (IP) address), some other globally
 22 unique identifier, a locally unique identifier (e.g., a numbering scheme proprietary
 23 or local to co-location facility 208), etc.

24 Filters 114 can fully restrict access to a node (e.g., no data can be received
 25 from or sent to the node), or partially restrict access to a node. Partial access

1 restriction can take different forms. For example, a node may be restricted so that
2 data can be received from the node but not sent to the node (or vice versa). By
3 way of another example, a node may be restricted so that only certain types of data
4 (e.g., communications in accordance with certain protocols, such as HTTP) can be
5 received from and/or sent to the node. Filtering based on particular types of data
6 can be implemented in different manners, such as by communicating data in
7 packets with header information that indicate the type of data included in the
8 packet.

9 Filters 114 can be added by application operations management console
10 242 or cluster operations management console 240 of Fig. 4. In the illustrated
11 example, filters added by cluster operations management console 240 (to establish
12 cluster boundaries) restrict full access to nodes (e.g., any access to another node
13 can be prevented) whereas filters added by application operations management
14 console 242 (to establish sub-boundaries within a cluster) can restrict either full
15 access to nodes or partial access.

16 Controller 254 also imposes some restrictions on what filters can be added
17 to filters 114. In the illustrated example, controller 254 allows cluster operations
18 management console 240 to add any filters it desires (which will define the
19 boundaries of the cluster). However, controller 254 restricts application operations
20 management console 242 to adding only filters that are at least as restrictive as
21 those added by console 240. If console 242 attempts to add a filter that is less
22 restrictive than those added by console 240 (in which case the sub-boundary may
23 extend beyond the cluster boundaries), controller 254 refuses to add the filter (or
24 alternatively may modify the filter so that it is not less restrictive). By imposing
25 such a restriction, controller 254 can ensure that the sub-boundaries established at

DHCP module 260 implements the Distributed Host Control Protocol, allowing BMonitor 250 (and thus node 210) to obtain an IP address from a DHCP server (e.g., cluster operations management console 240 of Fig. 4). During an initialization process for node 210, DHCP module 260 requests an IP address from the DHCP server, which in turn provides the IP address to module 260. Additional information regarding DHCP is available from Microsoft Corporation of Redmond, Washington.

Software engines 252 include any of a wide variety of conventional software components. Examples of engines 252 include an operating system (e.g., Windows NT®), a load balancing server component (e.g., to balance the processing load of multiple nodes 210), a caching server component (e.g., to cache data and/or instructions from another node 210 or received via the Internet), a storage manager component (e.g., to manage storage of data from nodes 210 received via the Internet), etc. In one implementation, each of the engines 252 is a protocol-based engine, communicating with BMonitor 250 and other engines 252 via messages and/or function calls without requiring the engines 252 and BMonitor 250 to be written using the same programming language.

Controller 254 may optionally be further responsible for controlling the execution of engines 252. This control can take different forms, including beginning execution of an engine 252, terminating execution of an engine 252, re-loading an image of an engine 252 from a storage device, etc. Controller 254 receives instructions from application operations management console 242 of Fig. 4 regarding which of these control actions to take and when to take them. Thus,

the control of engines 252 is actually managed by the remote application operations management console 242, not locally at co-location facility 104. Controller 254 also provides an interface via which application operations management console 242 can identify filters to add (and/or remove) from filter set 258.

Controller 254 also includes an interface via which cluster operations management console 240 of Fig. 4 can communicate commands to controller 254. Different types of hardware operation oriented commands can be communicated to controller 254 by cluster operations management console 240, such as re-booting the node, shutting down the node, placing the node in a low-power state (e.g., in a suspend or standby state), etc.

Controller 254 further optionally provides encryption support for BMonitor 250, allowing data to be stored securely on mass storage device 262 (e.g., a magnetic disk, an optical disk, etc.) and secure communications to occur between node 210 and an operations management console (e.g., console 240 or 242 of Fig. 4). Controller 254 maintains multiple encryption keys 258, which can include a variety of different keys such as symmetric keys (secret keys used in secret key cryptography), public/private key pairs (for public key cryptography), etc. to be used in encrypting and/or decrypting data.

Fig. 6 is a flowchart illustrating an exemplary process for making modifications to restrictions in a network mediator in accordance with certain embodiments of the invention. The process of Fig. 6 is performed by a network mediator (such as mediator 108 or 110) and may be implemented in software.

Initially, the network mediator authenticates a remote manager (act 280). If the remote device attempting to be authenticated cannot be authenticated, then no

modifications to the restrictions on the network mediator are allowed (acts 282 and 284). However, if the remote device can be authenticated, then a request to modify restrictions can be received by the network mediator (acts 282 and 286). A request to modify a restriction can be a request to modify a filter and/or virtualization data, such as to add a filter or virtualization data, remove a filter or virtualization data, modify the parameters of a filter or virtualization data, etc.

In response to the request to modify a restriction, the network mediator checks whether the requested modification could result in a violation of a higher-level restriction (act 288). If a violation could result, then the requested modification is denied (act 290). However, if a violation would not result, then the requested modification is made (act 292).

Fig. 7 is a flowchart illustrating an exemplary process for imposing restrictions on communications of a computing device in accordance with certain embodiments of the invention. The process of Fig. 7 is performed by a network mediator (such as mediator 108 or 110) and may be implemented in software.

Initially, a data packet is received at the network mediator (act 310). The received packet could be from the computing device corresponding to the network mediator and targeting some other computing device, or alternatively from some other computing device targeting the computing device corresponding to the network mediator. Alternatively, data packets in only one direction (e.g., from the computing device corresponding to the network mediator out to some other computing device, or from some other computing device targeting the computing device corresponding to the network mediator). Once received, the network mediator determines whether the filters indicate it is okay to pass the packet through to its target (act 312). If the filters indicate it is not okay to pass the

1 packet, then the packet is dropped (act 314). The source of the packet may be
2 informed that the packet was dropped, or alternatively may not be so informed.

3 However, if the filters indicate it is okay to pass the packet through to its
4 target, then the network mediator checks whether virtualized addresses are being
5 used, either for the network mediator as a whole, or alternatively on an individual
6 per-address (or per-filter) basis (act 316). If there is no virtualized address for this
7 packet, then the data packet is passed to the addressed device (act 318). However,
8 if there is a virtualized address for this packet, then the network mediator replaces
9 the address in the data packet with the mapped address (act 320). This
10 replacement is either replacing a virtual address with a network address (in the
11 case of a data packet being sent from the computing device corresponding to the
12 network mediator), or replacing a network address with a virtual address (in the
13 case of a data packet targeting the computing device corresponding to the network
14 mediator).

15 In the discussion herein, embodiments of the invention are described in the
16 general context of computer-executable instructions, such as program modules,
17 being executed by one or more conventional personal computers (e.g., computing
18 devices 102 and 104 of Fig. 1). Generally, program modules include routines,
19 programs, objects, components, data structures, etc. that perform particular tasks
20 or implement particular abstract data types. Moreover, those skilled in the art will
21 appreciate that various embodiments of the invention may be practiced with other
22 computer system configurations, including hand-held devices, gaming consoles,
23 Internet appliances, multiprocessor systems, microprocessor-based or
24 programmable consumer electronics, network PCs, minicomputers, mainframe
25

computers, and the like. In a distributed computer environment, program modules may be located in both local and remote memory storage devices.

Alternatively, embodiments of the invention can be implemented in hardware or a combination of hardware, software, and/or firmware. For example, all or part of the invention can be implemented in one or more application specific integrated circuits (ASICs) or programmable logic devices (PLDs).

Fig. 8 shows a general example of a computer 342 that can be used in accordance with certain embodiments of the invention. Computer 342 is shown as an example of a computer that can perform the functions of a computing device 102 or 104 of Fig 1, or node 210 of Fig. 4 or 5.

Computer 342 includes one or more processors or processing units 344, a system memory 346, and a bus 348 that couples various system components including the system memory 346 to processors 344. The bus 348 represents one or more of any of several types of bus structures, including a memory bus or memory controller, a peripheral bus, an accelerated graphics port, and a processor or local bus using any of a variety of bus architectures. The system memory includes read only memory (ROM) 350 and random access memory (RAM) 352. A basic input/output system (BIOS) 354, containing the basic routines that help to transfer information between elements within computer 342, such as during start-up, is stored in ROM 350.

Computer 342 further includes a hard disk drive 356 for reading from and writing to a hard disk, not shown, connected to bus 348 via a hard disk driver interface 357 (e.g., a SCSI, ATA, or other type of interface); a magnetic disk drive 358 for reading from and writing to a removable magnetic disk 360, connected to bus 348 via a magnetic disk drive interface 361; and an optical disk drive 362 for

A number of program modules may be stored on the hard disk, magnetic disk 360, optical disk 364, ROM 350, or RAM 352, including an operating system 370, one or more application programs 372, other program modules 374, and program data 376. A user may enter commands and information into computer 342 through input devices such as keyboard 378 and pointing device 380. Other input devices (not shown) may include a microphone, joystick, game pad, satellite dish, scanner, or the like. These and other input devices are connected to the processing unit 344 through an interface 368 that is coupled to the system bus. A monitor 384 or other type of display device is also connected to the system bus 348 via an interface, such as a video adapter 386. In addition to the monitor, personal computers typically include other peripheral output devices (not shown) such as speakers and printers.

Computer 342 optionally operates in a networked environment using logical connections to one or more remote computers, such as a remote computer

388. The remote computer 388 may be another personal computer, a server, a router, a network PC, a peer device or other common network node, and typically includes many or all of the elements described above relative to computer 342, although only a memory storage device 390 has been illustrated in Fig. 8. The logical connections depicted in Fig. 8 include a local area network (LAN) 392 and a wide area network (WAN) 394. Such networking environments are commonplace in offices, enterprise-wide computer networks, intranets, and the Internet. In the described embodiment of the invention, remote computer 388 executes an Internet Web browser program (which may optionally be integrated into the operating system 370) such as the "Internet Explorer" Web browser manufactured and distributed by Microsoft Corporation of Redmond, Washington.

When used in a LAN networking environment, computer 342 is connected to the local network 392 through a network interface or adapter 396. When used in a WAN networking environment, computer 342 typically includes a modem 398 or other component for establishing communications over the wide area network 394, such as the Internet. The modem 398, which may be internal or external, is connected to the system bus 348 via an interface (e.g., a serial port interface 368). In a networked environment, program modules depicted relative to the personal computer 342, or portions thereof, may be stored in the remote memory storage device. It is to be appreciated that the network connections shown are exemplary and other means of establishing a communications link between the computers may be used.

Computer 342 also includes a broadcast tuner 400. Broadcast tuner 400 receives broadcast signals either directly (e.g., analog or digital cable

1 transmissions fed directly into tuner 400) or via a reception device (e.g., via an
2 antenna or satellite dish).

3 Generally, the data processors of computer 342 are programmed by means
4 of instructions stored at different times in the various computer-readable storage
5 media of the computer. Programs and operating systems are typically distributed,
6 for example, on floppy disks or CD-ROMs. From there, they are installed or
7 loaded into the secondary memory of a computer. At execution, they are loaded at
8 least partially into the computer's primary electronic memory. The invention
9 described herein includes these and other various types of computer-readable
10 storage media when such media contain instructions or programs for implementing
11 the acts described herein in conjunction with a microprocessor or other data
12 processor. The invention also includes the computer itself when programmed
13 according to the methods and techniques described herein. Furthermore, certain
14 sub-components of the computer may be programmed to perform the functions
15 and steps described herein. The invention includes such sub-components when
16 they are programmed as described. In addition, the invention described herein
17 includes data structures as embodied on various types of memory media.

18 For purposes of illustration, programs and other executable program
19 components such as the operating system are illustrated herein as discrete blocks,
20 although it is recognized that such programs and components reside at various
21 times in different storage components of the computer, and are executed by the
22 data processor(s) of the computer.

1 **Conclusion**

2 Although the description above uses language that is specific to structural
3 features and/or methodological acts, it is to be understood that the invention
4 defined in the appended claims is not limited to the specific features or acts
5 described. Rather, the specific features and acts are disclosed as exemplary forms
6 of implementing the invention.

7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25